

Федеральное государственное автономное образовательное
учреждение высшего образования
«СИБИРСКИЙ ФЕДЕРАЛЬНЫЙ УНИВЕРСИТЕТ»
Институт фундаментальной биологии и биотехнологий
Кафедра биофизики

УТВЕРЖДАЮ

Заведующий кафедрой

_____ / Кратасюк В.А.

«___» _____ 2016 г.

БАКАЛАВРСКАЯ РАБОТА

06.03.01. Биология

**ИЗМЕНЧИВОСТЬ ГЕНОМОВ ПО СТЕПЕНИ НАРУШЕНИЯ ВТОРОГО
ПРАВИЛА ЧАРГАФФА**

Руководитель _____

д.ф.-м.н. М. Г. Садовский

Выпускник _____

Я. В. Петрова

Красноярск 2016

СОДЕРЖАНИЕ

ВВЕДЕНИЕ.....	3
1. ОБЗОР ЛИТЕРАТУРЫ.....	4
1.1. Строение ДНК.....	4
1.2. Длина генома.....	4
1.3. Частотный словарь.....	6
1.4. Правила Чаргаффа.....	6
2. МАТЕРИАЛЫ И МЕТОДЫ ИССЛЕДОВАНИЯ.....	7
2.1. Материалы.....	7
2.2. Методы.....	20
3. РЕЗУЛЬТАТЫ.....	20
4. ВЫВОДЫ.....	54
СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ.....	55

ВВЕДЕНИЕ

Наука не стоит на месте, еще совсем недавно только начинали секвенировать геномы простейших организмов (*Haemophilus influenzae* 1995г.) [1], а на сегодняшний день уже очень много полностью аннотированных геномов различных организмов (включая человека), они нуждаются в обработке, но работать с ними неудобно из-за различной длины геномов.

Поэтому придумали простой способ сократить разницу между геномами посредством перехода от полного текста к его частотному словарю. Это простой, но продуктивный прием, позволяющий единообразно работать с текстами различной длины, сравнивать их, производить информационный анализ.

Цель моей работы заключалась в выявлении наличия зависимости между таксономическими признаками и степенью нарушения второго правила Чаргаффа в геноме.

Были поставлены следующие задачи:

1. Составить базу данных геномов бактерий, архей, эукариот.
2. Сравнить, совпадают ли разделения на группы по коэффициентам невязки с таксономией.

1. ОБЗОР ЛИТЕРАТУРЫ

1.1. Строение ДНК

С химической точки зрения ДНК — это длинная полимерная молекула, состоящая из повторяющихся блоков — нуклеотидов. Каждый нуклеотид состоит из азотистого основания, сахара (дезоксирибозы) и фосфатной группы. Связи между нуклеотидами в цепи образуются за счёт дезоксирибозы и фосфатной группы (фосфодиэфирные связи). В подавляющем большинстве случаев (кроме некоторых вирусов, содержащих одноцепочечную ДНК) макромолекула ДНК состоит из двух цепей, ориентированных азотистыми основаниями друг к другу. Эта двухцепочечная молекула спирализована. В целом структура молекулы ДНК получила название «двойной спирали».

В ДНК встречается четыре вида азотистых оснований (аденин, гуанин, тимин и цитозин). Азотистые основания одной из цепей соединены с азотистыми основаниями другой цепи водородными связями согласно принципу комплементарности: аденин соединяется только с тимином, гуанин — только с цитозином.

Последовательность нуклеотидов позволяет «кодировать» информацию о различных типах РНК, наиболее важными из которых являются информационные, или матричные (мРНК), рибосомальные (рРНК) и транспортные (тРНК). Все эти типы РНК синтезируются на матрице ДНК за счёт копирования последовательности ДНК в последовательность РНК, синтезируемой в процессе транскрипции, и принимают участие в биосинтезе белков (процессе трансляции). Помимо кодирующих последовательностей, ДНК клеток содержит последовательности, выполняющие регуляторные и структурные функции. Кроме того, в геноме эукариот часто встречаются участки, принадлежащие «генетическим паразитам», например, транспозонам [2].

1.2. Длина генома

Термин «геном» был предложен Г. Винклером в 1920 г. для описания совокупности генов, заключенных в гаплоидном наборе хромосом одного

биологического вида. Уже в то время подчеркивалось, что понятие генома в отличие от генотипа является характеристикой вида в целом, а не отдельной особи. Интенсивное исследование геномов в течение последних пятидесяти лет заметно изменило представления об их строении и функционировании. С установлением полной первичной структуры генома человека, а также целого ряда животных, растений и микроорганизмов человечество вступило в «постгеномную эру». Сегодня термином «геном» обозначают совокупность ДНК гаплоидного набора хромосом, которые заключены в отдельной клетке зародышевой линии многоклеточного организма[3].

По соотношению размера генома и числа генов геномы могут быть разделены на два чётко выделенных класса:

1. Небольшие компактные геномы размером, как правило, не более 10 млн пар оснований, со строгим соответствием между размером генома и числом генов. Такими геномами обладают все вирусы и прокариоты. У этих организмов плотность генов составляет от 0,5 до 2 генов на тысячу пар оснований, а между генами имеются очень короткие участки, занимающие 10-15 % длины генома. Межгенные участки в таких геномах состоят главным образом из регуляторных элементов. Помимо вирусов и прокариот к этому классу могут быть отнесены и геномы большинства одноклеточных эукариот, хотя их геномы демонстрируют несколько меньшую зависимость между размером генома и числом генов, а размер генома может достигать 20 млн пар оснований.
2. Обширные геномы размером более 100 млн пар оснований, у которых нет чёткой взаимосвязи между размером генома и числом генов. К этому классу относятся большие геномы многоклеточных эукариот и некоторых одноклеточных эукариот. В отличие от геномов первой группы большинство нуклеотидов в геномах этого класса относятся к последовательностям, которые не кодируют ни белков, ни РНК[4].

Для того, чтобы понять, насколько сильно может различаться длина генома у разных организмов, приведу пример самого маленького генома - бактерии *Carsonella ruddii*, длина генома составляет 159 662 пары оснований [5], самым длинным расшифрованным геномом обладает ладанная сосна *Pinus taeda*, длина генома составляет 22,18 миллиардов пар оснований [6].

1.3. Частотный словарь

Нуклеотидные последовательности мы рассматриваем как текст, т.е. как последовательность символов из четырехбуквенного алфавита А,С,Г,Т. Длина генетического текста — количество N нуклеотидов в нуклеотидной последовательности.

Предполагается, что клеточные механизмы, связанные со считыванием наследственной информации, оперируют с отдельными, весьма малыми, фрагментами ДНК. В соответствии с этим перейдем от рассмотрения молекулы как целого к изучению совокупности ее фрагментов фиксированной длины (слов).

Список всех слов длины q , входящих в данный текст, называется q -носителем данного текста. Если каждому слову q -носителя сопоставить частоту его встречаемости в тексте, получим частотный словарь длины q [7]. Мы будем рассматривать словари при $q=2-6$ нуклеотидов.

Частоты слов были рассчитаны по следующей формуле:

$$f_{\omega} = \frac{n_{\omega}}{N} \quad (1)$$

где n_{ω} — число копий каждого слова, ω — слово, N — общее число слов [8].

1.4. Правила Чаргаффа

Правила Чаргаффа — система эмпирически выявленных правил, описывающих количественные соотношения между различными типами азотистых оснований в ДНК. Были сформулированы в результате работы группы биохимика Эрвина Чаргаффа в 1949-1951 гг.

До работ группы Чаргаффа господствовала так называемая «тетрануклеотидная» теория, согласно которой ДНК состоит из повторяющихся блоков по четыре разных азотистых основания (аденин, тимин, гуанин и цитозин). Чаргаффу и сотрудникам удалось разделить нуклеотиды ДНК при помощи бумажной хроматографии и определить точные количественные соотношения нуклеотидов разных типов. Они

значительно отличались от эквимольных, которых можно было бы ожидать, если бы все четыре основания были представлены в равных пропорциях.

В 1968 году Чарграфф выявил, что в каждой из нитей ДНК количество аденина приблизительно равно количеству тимина, а гуанина — цитозину: $A \sim T$, $G \sim C$. В 1990-х с развитием технологии секвенирования ДНК это правило было подтверждено.

Соотношения, выявленные Чарграффом для аденина (А), тимина (Т), гуанина (Г) и цитозина (Ц), оказались следующими:

1. Количество аденина равно количеству тимина, а гуанина — цитозину: $A=T$, $G=C$.
2. Количество пуринов равно количеству пиримидинов: $A+G=T+C$.
3. Количество оснований с аминогруппами в положении 6 равно количеству оснований с кетогруппами в положении 6: $A+C=G+T$.

Вместе с тем, соотношение $(A+T):(G+C)$ может быть различным у ДНК разных видов. У одних преобладают пары АТ, в других — ГЦ. Эти соотношения применялись только к одноцепочечной ДНК.

Правила Чаргаффа, наряду с данными рентгеноструктурного анализа, сыграли решающую роль в расшифровке структуры ДНК Дж. Уотсоном и Фрэнсисом Криком.

Позднее, когда было секвенировано больше геномов, стало понятно, что соотношение $(A+T):(G+C)$ выполняется и для обеих нитей ДНК одновременно. Это и было названо вторым правилом Чаргаффа[9,10].

Известно, что второе правило Чаргаффа нарушается у геномов разных организмов по-разному. Особенно часто нарушается оно в митохондриальных геномах [11].

2. МАТЕРИАЛЫ И МЕТОДЫ ИССЛЕДОВАНИЯ

2.1. Материалы

Ниже приведены три таблицы с перечислением геномов, для которых определялась невязка нарушения второго правила Чаргаффа.

Геномы взяты из базы данных EMBL - <http://www.embl.org/>.

Таблица 1. Геномы архей, всего 109.

Отдел А1 <i>Crenarchaeota</i>	Отдел А2 <i>Euryarchaeota</i>
<i>Acidianus hospitalis</i>	<i>Archaeoglobus fulgidus</i>
<i>Acidilobus saccharovorans</i>	<i>Archaeoglobus profundus</i>
<i>Aeropyrum camini</i>	<i>Archaeoglobus sulfaticallidus</i>
<i>Aeropyrum pernix</i>	<i>Archaeoglobus veneficus</i>
<i>Caldisphaera lagunensis</i>	<i>Methanomassiliicoccus intestinalis</i>
<i>Caldivirga maquilingensis</i>	<i>Methanomethylophilus alvus</i>
<i>Desulfurococcus fermentans</i>	<i>Methanoplasma termitum</i>
<i>Desulfurococcus mucosus</i>	<i>Ferroglobus placidus</i>
<i>Fervidicoccus fontis</i>	<i>Ferroplasma acidarmanus</i>
<i>Hyperthermus butylicus</i>	<i>Geoglobus acetivorans</i>
<i>Ignicoccus hospitalis</i>	<i>Halalkalicoccus jeotgali</i>
<i>Ignisphaera aggregans</i>	<i>Haloarcula hispanica</i>
<i>Metallosphaera cuprina</i>	<i>Halobacterium salinarum</i>
<i>Metallosphaera sedula</i>	<i>Haloferax mediterranei</i>
<i>Pyrobaculum arsenaticum</i>	<i>Haloferax volcanii</i>
<i>Pyrobaculum calidifontis</i>	<i>Halogeometricum borinquense</i>
<i>Pyrolobus fumarii</i>	<i>Halomicrobium mukohataei</i>
<i>Sulfolobus tokodaii</i>	<i>Halopiger xanaduensis</i>
<i>Thermogladius cellulolyticus</i>	<i>Haloquadratum walsbyi</i>
<i>Thermoproteus tenax</i>	<i>Halorhabdus tiamatea</i>
<i>Vulcanisaeta distributa</i>	<i>Halorubrum lacusprofundi</i>
<i>Vulcanisaeta moutnovskia</i>	<i>Halostagnicola larsenii</i>
	<i>Halovivax ruber</i>
Отдел А3 <i>Korarchaeota</i>	<i>Methanobacterium formicicum</i>
<i>Korarchaeum cryptofilum</i>	<i>Methanobacterium lacus</i>

	<i>Methanobacterium paludis</i>
Отдел A4 Nanoarchaeota	<i>Methanobrevibacter ruminantium</i>
<i>Nanoarchaeum equitans</i>	<i>Methanobrevibacter smithii</i>
	<i>Methanocaldococcus fervens</i>
Без отдела	<i>Methanocaldococcus infernus</i>
<i>Nitrosopumilus koreensis</i>	<i>Methanocaldococcus jannaschii</i>
<i>Nitrososphaera evergladensis</i>	<i>Methanocaldococcus vulcanius</i>
<i>Cenarchaeum symbiosum</i>	<i>Methanocella arvoryzae</i>
<i>Nitrosopumilus maritimus</i>	<i>Methanocella conradii</i>
	<i>Methanocella paludicola</i>
	<i>Methanococcoides burtonii</i>
	<i>Methanococcoides methylutens</i>
	<i>Methanococcus aeolicus</i>
	<i>Methanococcus vanniellii</i>
	<i>Methanococcus voltae</i>
	<i>Methanocorpusculum labreanum</i>
	<i>Methanoculleus bourgensis</i>
	<i>Methanoculleus marisnigri</i>
	<i>Methanohalobium evestigatum</i>
	<i>Methanohalophilus mahii</i>
	<i>Methanolacinia petrolearia</i>
	<i>Methanolobus psychrophilus</i>
	<i>Methanomethylovorans hollandica</i>
	<i>Methanopyrus kandleri</i>
	<i>Methanoregula boonei</i>
	<i>Methanoregula formicica</i>
	<i>Methanosaeta concilii</i>

	<i>Methanosaeta harundinacea</i>
	<i>Methanosalsum zhilinae</i>
	<i>Methanosarcina acetivorans</i>
	<i>Methanosarcina thermophila</i>
	<i>Methanosarcina vacuolata</i>
	<i>Methanosphaera stadtmanae</i>
	<i>Methanosphaerula palustris</i>
	<i>Methanospirillum hungatei</i>
	<i>Methanothermobacter marburgensis</i>
	<i>Methanothermobacter thermautotrophicus</i>
	<i>Methanothermococcus okinawensis</i>
	<i>Methanothermus fervidus</i>
	<i>Methanotorris igneus</i>
	<i>Natrialba magadii</i>
	<i>Natrinema pellirubrum</i>
	<i>Natronobacterium gregoryi</i>
	<i>Natronococcus occultus</i>
	<i>Natronomonas moolapensis</i>
	<i>Natronomonas pharaonis</i>
	<i>Picrophilus torridus</i>
	<i>Pyrococcus abyssi</i>
	<i>Pyrococcus furiosus</i>
	<i>Pyrococcus horikoshii</i>
	<i>Thermococcus barophilus</i>
	<i>Thermococcus kodakarensis</i>
	<i>Thermococcus litoralis</i>
	<i>Thermococcus nautili</i>

	<i>Thermoplasma acidophilum</i>
	<i>Thermoplasma volcanium</i>

Таблица 2. Геномы бактерий, всего 240.

Без отдела	Отдел B12 <i>Proteobacteria</i>
<i>Aminobacterium colombiense</i>	<i>Acetobacter pasteurianus</i>
<i>Anaerobaculum mobile</i>	<i>Achromobacter xylosoxidans</i>
<i>Candidatus Koribacter versatilis</i>	<i>Acidovorax avenae</i>
<i>Marivirga tractuosa</i>	<i>Acidovorax citrulli</i>
<i>Oceanithermus profundus</i>	<i>Acinetobacter calcoaceticus</i>
<i>Pseudopedobacter saltans</i>	<i>Acinetobacter oleivorans</i>
	<i>Actinobacillus pleuropneumoniae</i>
Отдел B1 <i>Aquificae</i>	<i>Actinobacillus succinogenes</i>
<i>Aquifex aeolicus</i>	<i>Actinobacillus suis</i>
	<i>Advenella kashmirensis</i>
Отдел B2 <i>Thermotogae</i>	<i>Advenella mimigardefordensis</i>
<i>Fervidobacterium pennivorans</i>	<i>Aeromonas media</i>
<i>Pseudothermotoga thermarum</i>	<i>Aeromonas veronii</i>
	<i>Aggregatibacter aphrophilus</i>
Отдел B4 <i>Deinococcus-Thermus</i>	<i>Alcanivorax borkumensis</i>
<i>Deinococcus maricopensis</i>	<i>Alcanivorax dieselolei</i>
<i>Marinithermus hydrothermalis</i>	<i>Alcanivorax pacificus</i>
<i>Meiothermus ruber</i>	<i>Alicyclophilus denitrificans</i>
<i>Meiothermus silvanus</i>	<i>Allochromatium vinosum</i>
	<i>Arcobacter nitrofigilis</i>
Отдел B6 <i>Chloroflexi</i>	<i>Arthrobacter chlorophenolicus</i>

<i>Chloroflexus aurantiacus</i>	<i>Azorhizobium caulinodans</i>
<i>Herpetosiphon aurantiacus</i>	<i>Azotobacter chroococcum</i>
<i>Roseiflexus castenholzii</i>	<i>Bacillus anthracis</i>
	<i>Bacillus cellulosilyticus</i>
Отдел B8 <i>Nitrospira</i>	<i>Bacillus cereus</i>
<i>Desulfovibrio salexigens</i>	<i>Bacillus methanolicus</i>
<i>Leptospirillum ferriphilum</i>	<i>Bacillus subtilis</i>
	<i>Bacillus thuringiensis</i>
Отдел B9 <i>Deferribacteres</i>	<i>Bacteriovorax marinus</i>
<i>Denitrovibrio acetiphilus</i>	<i>Bdellovibrio exovorus</i>
<i>Flexistipes sinusarabici</i>	<i>Bordetella avium</i>
	<i>Bordetella holmesii</i>
Отдел B10 <i>Cyanobacteria</i>	<i>Campylobacter coli</i>
<i>Acaryochloris marina</i>	<i>Campylobacter curvus</i>
<i>Anabaena cylindrica</i>	<i>Campylobacter volucris</i>
<i>Chroococcidiopsis thermalis</i>	<i>Candidatus Kinetoplastibacterium galatii</i>
<i>Dactylococcopsis salina</i>	<i>Candidatus Portiera aleyrodidarum</i>
	<i>Candidatus Symbiobacter mobilis</i>
Отдел B11 <i>Chlorobi</i>	<i>Chromohalobacter salexigens</i>
<i>Pelodictyon luteolum</i>	<i>Citrobacter freundii</i>
	<i>Coxiella burnetii</i>
Отдел B13 <i>Firmicutes</i>	<i>Desulfobacca acetoxidans</i>
<i>Acetobacterium woodii</i>	<i>Desulfohalobium retbaense</i>
<i>Acetohalobium arabaticum</i>	<i>Desulfurobacterium thermolithotrophum</i>
<i>Acholeplasma laidlawii</i>	<i>Erwinia pyrifoliae</i>
<i>Acholeplasma oculi</i>	<i>Escherichia coli</i>
<i>Acidaminococcus fermentans</i>	<i>Ferrimonas balearica</i>

<i>Ammonifex degensii</i>	<i>Geobacter metallireducens</i>
<i>Anaerococcus prevotii</i>	<i>Gluconobacter oxydans</i>
<i>Anoxybacillus flavithermus</i>	<i>Haemophilus ducreyi</i>
<i>Bacteroides salanitronis</i>	<i>Haemophilus influenzae</i>
<i>Bacteroides vulgatus</i>	<i>Haliangium ochraceum</i>
<i>Brevibacillus brevis</i>	<i>Halomonas elongata</i>
<i>Brevibacillus laterosporus</i>	<i>Helicobacter cetorum</i>
<i>Caldicellulosiruptor hydrothermalis</i>	<i>Helicobacter mustelae</i>
<i>Caldicellulosiruptor kristjanssonii</i>	<i>Helicobacter pylori</i>
<i>Clostridium botulinum</i>	<i>Hipaea maritima</i>
<i>Dehalobacter restrictus</i>	<i>Kangiella koreensis</i>
<i>Desulfitobacterium hafniense</i>	<i>Legionella pneumophila</i>
<i>Desulfosporosinus meridiei</i>	<i>Leptothrix cholodnii</i>
<i>Desulfotomaculum acetoxidans</i>	<i>Mannheimia haemolytica</i>
<i>Desulfotomaculum gibsoniae</i>	<i>Myxococcus stipitatus</i>
<i>Desulfotomaculum ruminis</i>	<i>Neisseria gonorrhoeae</i>
<i>Exiguobacterium sibiricum</i>	<i>Neisseria meningitidis</i>
<i>Geobacillus kaustophilus</i>	<i>Nitratifractor salsuginis</i>
<i>Geobacillus thermodenitrificans</i>	<i>Nitrosomonas eutropha</i>
<i>Halanaerobium praevalens</i>	<i>Novosphingobium aromaticivorans</i>
<i>Halobacillus halophilus</i>	<i>Paracoccus aminophilus</i>
<i>Halobacteroides halobius</i>	<i>Pectobacterium wasabiae</i>
<i>Kyrpidia tusciae</i>	<i>Phaeobacter gallaeciensis</i>
<i>Lactobacillus casei</i>	<i>Phaeobacter inhibens</i>
<i>Lactobacillus helveticus</i>	<i>Pseudomonas aeruginosa</i>
<i>Lactobacillus reuteri</i>	<i>Pseudomonas balearica</i>
<i>Lactobacillus rhamnosus</i>	<i>Pseudomonas fluorescens</i>

<i>Lactococcus lactis</i>	<i>Pseudomonas mosselii</i>
<i>Listeria ivanovii</i>	<i>Pseudomonas putida</i>
<i>Listeria monocytogenes</i>	<i>Ralstonia solanacearum</i>
<i>Mycoplasma gallisepticum</i>	<i>Rhizobium etli</i>
<i>Mycoplasma hyorhinis</i>	<i>Roseibacterium elongatum</i>
<i>Mycoplasma pneumoniae</i>	<i>Salmonella enterica</i>
<i>Paenibacillus borealis</i>	<i>Starkeya novella</i>
<i>Paenibacillus durus</i>	<i>Xylella fastidiosa</i>
<i>Paenibacillus mucilaginosus</i>	<i>Yersinia pestis</i>
<i>Paenibacillus polymyxa</i>	<i>Yersinia pseudotuberculosis</i>
<i>Paenibacillus sabinae</i>	<i>Yersinia ruckeri</i>
<i>Paenibacillus terrae</i>	
<i>Peptoclostridium difficile</i>	Отдел B14 Bacteroidetes
<i>Roseburia hominis</i>	<i>Aequorivita sublithicola</i>
<i>Staphylococcus aureus</i>	<i>Alistipes finegoldii</i>
<i>Staphylococcus pseudintermedius</i>	<i>Barnesiella viscericola</i>
<i>Streptococcus agalactiae</i>	<i>Capnocytophaga canimorsus</i>
<i>Streptococcus mutans</i>	<i>Cellulophaga baltica</i>
<i>Streptococcus pneumoniae</i>	<i>Cellulophaga lytica</i>
<i>Streptococcus pseudopneumoniae</i>	<i>Dyadobacter fermentans</i>
<i>Streptococcus pyogenes</i>	<i>Emticicia oligotrophica</i>
<i>Veillonella parvula</i>	<i>Flexibacter litoralis</i>
<i>Weissella ceti</i>	<i>Fluviicola taffensis</i>
	<i>Haliscomenobacter hydrossis</i>
	<i>Muricauda ruestringensis</i>
Отдел B15 Planctomycetes	<i>Niabella soli</i>
<i>Pirellula staleyi</i>	<i>Odoribacter splanchnicus</i>

<i>Planctopirus limnophilus</i>	<i>Ornithobacterium rhinotracheale</i>
<i>Rubinisphaera brasiliensis</i>	<i>Owenweeksia hongkongensis</i>
	<i>Pedobacter heparinus</i>
Отдел B16 <i>Chlamydiae</i>	<i>Pelobacter carbinolicus</i>
<i>Chlamydia avium</i>	<i>Pelobacter propionicus</i>
<i>Chlamydia pecorum</i>	<i>Porphyromonas asaccharolytica</i>
	<i>Rhodothermus marinus</i>
Отдел B17 <i>Spirochaetes</i>	<i>Riemerella anatipestifer</i>
<i>Borrelia anserina</i>	<i>Salinibacter ruber</i>
<i>Borrelia bissettii</i>	<i>Weeksella virosa</i>
<i>Borrelia coriaceae</i>	
<i>Borrelia hermsii</i>	Отдел B24 <i>Actinobacteria</i>
<i>Borrelia recurrentis</i>	<i>Acidothermus cellulolyticus</i>
<i>Brachyspira hyodysenteriae</i>	<i>Actinoplanes missouriensis</i>
<i>Treponema brennaborensense</i>	<i>Actinosynnema mirum</i>
<i>Treponema caldaria</i>	<i>Adlercreutzia equolifaciens</i>
	<i>Amycolatopsis methanolica</i>
Отдел B19 <i>Fusobacteria</i>	<i>Amycolatopsis orientalis</i>
<i>Ilyobacter polytropus</i>	<i>Atopobium parvulum</i>
	<i>Bifidobacterium animalis</i>
Отдел B21 <i>Verrucomicrobia</i>	<i>Bifidobacterium bifidum</i>
<i>Akkermansia muciniphila</i>	<i>Blastococcus saxobsidens</i>
<i>Opitutus terrae</i>	<i>Brachybacterium faecium</i>
	<i>Corynebacterium glyciniphilum</i>
Отдел B22 <i>Dictyoglomi</i>	<i>Corynebacterium jeikeium</i>
<i>Dictyoglomus turgidum</i>	<i>Eggerthella lenta</i>
	<i>Geodermatophilus obscurus</i>

	<i>Intrasporangium calvum</i>
	<i>Isopterocola variabilis</i>
	<i>Jonesia denitrificans</i>
	<i>Kocuria rhizophila</i>
	<i>Kribbella flavida</i>
	<i>Kutzneria albida</i>
	<i>Kytococcus sedentarius</i>
	<i>Mycobacterium tuberculosis</i>
	<i>Nakamurella multipartita</i>
	<i>Nocardiopsis alba</i>
	<i>Parascardovia denticolens</i>
	<i>Rubrobacter xylanophilus</i>
	<i>Sanguibacter keddieii</i>
	<i>Streptomyces avermitilis</i>
	<i>Streptomyces cattleya</i>
	<i>Streptomyces fulvissimus</i>

Таблица 3. Геномы эукариот, всего 100.

<i>Anopheles gambiae</i> str. PEST chromosome X	<i>Kazachstania naganishii</i> CBS 8797 chromosome 4
<i>Arabidopsis thaliana</i> chromosome 4, long arm	<i>Kluyveromyces lactis</i> NRRL Y-1140 chromosome D
<i>Arabidopsis thaliana</i> chromosome 4, short arm	<i>Kluyveromyces marxianus</i> DMKU3-1042 DNA, chromosome 4
<i>Ashbya gossypii</i> ATCC 10895 chromosome IV	<i>Komagataella pastoris</i> CBS 7435 chromosome 4, complete replicon
<i>Aspergillus fumigatus</i> Af293 chromosome 4	<i>Lachancea kluyveri</i> NRRL Y-12651 chromosome D

<i>Aspergillus nidulans</i> FGSC A4 chromosome IV	<i>Lachancea thermotolerans</i> CBS 6340 chromosome D
<i>Aspergillus niger</i> supercontig Sc2, chromosome map 4R	<i>Leishmania braziliensis</i> chromosome 4
<i>Aspergillus niger</i> supercontig Sc7, chromosome map 4L	<i>Leishmania donovani</i> BPK282A1 chromosome 4
<i>Aspergillus niger</i> supercontig Sc19, chromosome map 4ER	<i>Leishmania major</i> strain Friedlin, chromosome 4
<i>Bos taurus</i> chromosome 4	<i>Leishmania mexicana</i> , chromosome 4
<i>Bos taurus</i> chromosome X	<i>Macaca fascicularis</i> chromosome 4
<i>Brachypodium distachyon</i> strain Bd21 chromosome 4	<i>Macaca fascicularis</i> chromosome X
<i>Caenorhabditis briggsae</i> AF16 draft chromosome, chrIV	<i>Macaca mulatta</i> chromosome 4
<i>Caenorhabditis briggsae</i> AF16 draft chromosome, chrX	<i>Macaca mulatta</i> chromosome X
<i>Caenorhabditis elegans</i> chromosome IV	<i>Medicago truncatula</i> chromosome 4
<i>Caenorhabditis elegans</i> chromosome X	<i>Nasonia vitripennis</i> chromosome 4
<i>Callithrix jacchus</i> chromosome 4	<i>Naumovozyma castellii</i> CBS 4309, chromosome 4
<i>Callithrix jacchus</i> chromosome X	<i>Naumovozyma dairenensis</i> CBS 421, chromosome 4
<i>Candida albicans</i> WO-1 chromosome 4 genomic scaffold	<i>Neospora caninum</i> Liverpool, chromosome IV
<i>Candida dubliniensis</i> CD36 chromosome 4	<i>Ornithorhynchus anatinus</i> chromosome 4
<i>Candida glabrata</i> CBS 138 chromosome D	<i>Oryctolagus cuniculus</i> chromosome 4
<i>Canis lupus familiaris</i> chromosome 4	<i>Oryza sativa</i> Indica Group chromosome 4
<i>Cryptococcus neoformans</i> var. <i>neoformans</i> B-3501A chromosome 4	<i>Oryzias latipes</i> chromosome 4
<i>Cryptococcus neoformans</i> var. <i>neoformans</i> JEC21 chromosome 4	<i>Ostreococcus lucimarinus</i> CCE9901 chromosome 4

<i>Cryptosporidium parvum</i> Iowa II chromosome 4	<i>Ovis aries</i> chromosome 4
<i>Cucumis melo</i> genomic chromosome, chr_4	<i>Pan troglodytes</i> chromosome 4
<i>Cyanidioschyzon merolae</i> strain 10D chromosome 4	<i>Phaeodactylum tricornutum</i> chromosome 4
<i>Danio rerio</i> chromosome 4	<i>Plasmodium falciparum</i> 3D7 chromosome 4
<i>Debaryomyces hansenii</i> CBS767 chromosome D	<i>Plasmodium knowlesi</i> strain H chromosome 4
<i>Dictyostelium discoideum</i> AX4 chromosome 4	<i>Plasmodium vivax</i> chromosome 4
<i>Drosophila melanogaster</i> chromosome 4	<i>Pongo abelii</i> chromosome 4
<i>Drosophila melanogaster</i> chromosome X	<i>Rattus norvegicus</i> chromosome 4
<i>Drosophila simulans</i> chromosome 4	<i>Saccharomyces cerevisiae</i> chromosome IV
<i>Drosophila simulans</i> chromosome X	<i>Scheffersomyces stipitis</i> CBS 6054 chromosome 4
<i>Drosophila yakuba</i> strain Tai18E2 chromosome 4	<i>Schistosoma mansoni</i> strain Puerto Rico chromosome 4
<i>Drosophila yakuba</i> strain Tai18E2 chromosome X	<i>Solanum lycopersicum</i> chromosome ch04
<i>Encephalitozoon cuniculi</i> GB-M1 chromosome IV	<i>Solanum pennellii</i> chromosome ch04
<i>Encephalitozoon intestinalis</i> ATCC 50506 chromosome IV	<i>Sorghum bicolor</i> chromosome 4
<i>Equus caballus</i> chromosome 4	<i>Sus scrofa</i> chromosome 4
<i>Equus caballus</i> chromosome X	<i>Taeniopygia guttata</i> chromosome 4
<i>Eremothecium cymbalariae</i> DBVPG#7215 chromosome 4	<i>Tetrapisispora blattae</i> CBS 6284 chromosome 4
<i>Fusarium graminearum</i> chromosome 4	<i>Tetrapisispora phaffii</i> CBS 4417 chromosome 4
<i>Fusarium oxysporum</i> f. sp. <i>lycopersici</i> 4287 chromosome 4	<i>Thalassiosira pseudonana</i> CCMP1335 chromosome 4

<i>Fusarium verticillioides</i> 7600 chromosome 4	<i>Theileria annulata</i> chromosome 4
<i>Gallus gallus</i> chromosome 4	<i>Theileria orientalis</i> strain Shintoku DNA, chromosome 4
<i>Glycine max</i> chromosome 4	<i>Thielavia terrestris</i> NRRL 8126 chromosome 4
<i>Gorilla gorilla gorilla</i> assembly, supercontig chr4	<i>Torulaspora delbrueckii</i> CBS 1146 chromosome 4
<i>Homo sapiens</i> chromosome 4	<i>Trypanosoma brucei gambiense</i> DAL972 chromosome 4
<i>Homo sapiens</i> chromosome X	<i>Yarrowia lipolytica</i> CLIB122 chromosome D
<i>Kazachstania africana</i> CBS 2517 chromosome 4	<i>Zea mays</i> chromosome 4
	<i>Zygosaccharomyces rouxii</i> strain CBS732 chromosome D

2.2. Методы

Коэффициент невязки был найден с помощью специального программного обеспечения по следующей формуле:

$$M_q = \frac{\sqrt{\sum_{\omega} (f_{\omega} - f_{\bar{\omega}})^2} * 1}{\|\Omega\|} \quad (2)$$

где $\Omega = \frac{4^q}{2}$; f_{ω} – частота встречаемости слова; $f_{\bar{\omega}}$ частота встречаемости комплиментарного слова в другой нити ДНК; ω – слово; q – длина слова.

Невязка измеряет меру отклонения внутри словаря, очень близка к Евклидовой метрике, но не является ею.

Если бы частоты слов были равными, то невязка была бы равна нулю и второе правило Чаргаффа выполнялось с абсолютной точностью [12].

3. РЕЗУЛЬТАТЫ

Ниже приведены графики, на которых сравниваются невязки разных видов организмов: архей, бактерий, эукариот.

Сравнивались значения невязок разных отделов при разной толщине словаря.

На графиках видно, что невязки родственных видов принимают близкие значения.



Рисунок 1 – Сравнение коэффициентов невязки для слов, длиной 1, разных геномов архей.



Рисунок 2 – Сравнение коэффициентов невязки для слов, длиной 2, разных геномов архей.



Рисунок 3 – Сравнение коэффициентов невязки для слов, длиной 3, разных геномов архей.



Рисунок 4 – Сравнение коэффициентов невязки для слов, длиной 4, разных геномов архей.



Рисунок 5 – Сравнение коэффициентов невязки для слов, длиной 5, разных геномов архей.



Рисунок 6 – Сравнение коэффициентов невязки для слов, длиной 6, разных геномов архей.



Рисунок 7 – Сравнение коэффициентов невязки для слов, длиной 1, разных геномов бактерий.



Рисунок 8 – Сравнение коэффициентов невязки для слов, длиной 2, разных геномов бактерий.



Рисунок 9 – Сравнение коэффициентов невязки для слов, длиной 3, разных геномов бактерий.



Рисунок 10 – Сравнение коэффициентов невязки для слов, длиной 4, разных геномов бактерий.



Рисунок 11 – Сравнение коэффициентов невязки для слов, длиной 5, разных геномов бактерий.



Рисунок 12 – Сравнение коэффициентов невязки для слов, длиной 6, разных геномов бактерий.



Рисунок 13 – Сравнение коэффициентов невязки для слов, длиной 1, разных геномов эукариот.



Рисунок 14 – Сравнение коэффициентов невязки для слов, длиной 2, разных геномов эукариот.



Рисунок 15 – Сравнение коэффициентов невязки для слов, длиной 3, разных геномов эукариот.



Рисунок 16 – Сравнение коэффициентов невязки для слов, длиной 4, разных геномов эукариот.



Рисунок 17 – Сравнение коэффициентов невязки для слов, длиной 5, разных геномов эукариот.



Рисунок 18 – Сравнение коэффициентов невязки для слов, длиной 6, разных геномов эукариот.

Ниже приведены графики, на которых сравниваются средние значения невязок разных видов организмов: архей, бактерий, эукариот; отделы разделены на порядки по таксономии.

На графиках видно, что значение невязки падает экспоненциально с ростом толщины словаря, это верно для всех организмов: архей, бактерий и эукариот.



Рисунок 19 – Среднее значение невязки отдела *Crenarchaeota*, разбитого на порядки.



Рисунок 20 – Среднее значение невязки отдела *Euryarchaeota*, разбитого на порядки.



Рисунок 21 – Значение невязки отдела *Korarchaeota*, вид *Korarchaeum cryptofilum*.



Рисунок 22 – Значение невязки отдела *Nanoarchaeota*, вид *Nanoarchaeum equitans*.



Рисунок 23 – Среднее значение невязки отдела *Thaumarchaeota*.



Рисунок 24 – Значение невязки отдела *Aquificae*, вид *Aquifex aeolicus*.



Рисунок 25 – Среднее значение невязки отдела *Thermotogae*, порядок *Thermotogales*.



Рисунок 26 – Среднее значение невязки отдела *Deinococcus-Thermus*.



Рисунок 27 – Среднее значение невязки отдела *Chloroflexi*.



Рисунок 28 – Среднее значение невязки отдела *Nitrospira*.



Рисунок 29 – Среднее значение невязки отдела *Deferribacteres*, порядок *Deferribacterales*.



Рисунок 30 – Среднее значение невязки отдела *Cyanobacteria*.



Рисунок 31 – Значение невязки отдела *Chlorobi*, вид *Pelodictyon luteolum*.



Рисунок 32 – Среднее значение невязки отдела *Proteobacteria*.



Рисунок 33 – Среднее значение невязки отдела *Firmicutes*.



Рисунок 34 – Среднее значение невязки отдела *Bacteroidetes*.



Рисунок 35 – Среднее значение невязки отдела *Planctomycetes*, порядок *Planctomycetales*.



Рисунок 36 – Среднее значение невязки отдела *Chlamydiae*, порядок *Chlamydiales*.



Рисунок 37 – Среднее значение невязки отдела *Spirochaetes*.



Рисунок 38 – Значение невязки отдела *Fusobacteria*, вид *Ilyobacter polytropus*.



Рисунок 39 – Среднее значение невязки отдела *Verrucomicrobia*.



Рисунок 40 – Значение невязки отдела *Dictyoglomi*, вид *Dictyoglomus turgidum*.



Рисунок 41 – Среднее значение невязки отдела *Actinobacteria*.



Рисунок 42 – Среднее значение невязки отдела *Synergistetes*.



Рисунок 43 – Среднее значение невязки отдела *Tenericutes*, класс *Mollicutes*.



Рисунок 44 – Значение невязки отдела *Acidobacteria*, вид *Koribacter versatilis*.



Рисунок 45 – Значение невязки отдела *Ascomycota*.



Рисунок 46 – Значение невязки отдела *Basidiomycota*, вид *Cryptococcus neoformans*.



Рисунок 47 – Значение невязки отдела *Microsporidia*, порядок *Microsporida*.



Рисунок 48 – Значение невязки отдела *Arthropoda*.



Рисунок 49 – Среднее значение невязки отдела *Chordata*.



Рисунок 50 – Среднее значение невязки отдела *Nematoda*, порядок *Rhabditida*.



Рисунок 51 – Значение невязки отдела *Platyhelminthes*, вид *Schistosoma mansoni*.



Рисунок 52 – Значение невязки отдела *Rhodophyta*, вид *Cyanidioschyzon merolae*.



Рисунок 53 – Значение невязки отдела *Amoebozoa*, вид *Dictyostelium discoideum*.



Рисунок 54 – Значение невязки отдела *Apicomplexa*.



Рисунок 55 – Значение невязки отдела *Euglenozoa*, порядок *Trypanosomatidae*.



Рисунок 56 – Значение невязок отдела *Heterokontophyta*.



Рисунок 57 – Значение невязки отдела *Chlorophyta*, вид *Ostreococcus lucimarinus*.



Рисунок 58 – Значение невязки отдела *Magnoliophyta*.

Так же была посчитана корреляция между значениями невязок разных порядков внутри отдела и между значениями невязок разных порядков разных отделов. Внутри отдела корреляция выше, чем между разными отделами.



Рисунок 59 – Корреляция между невязками внутри порядков разных отделов, и невязками между разными отделами архей.



Рисунок 60 – Корреляция между невязками внутри порядков разных отделов, и невязками между разными отделами бактерий.



Рисунок 61 – Корреляция между невязками внутри порядков разных отделов, и невязками между разными отделами архей.

4. ВЫВОДЫ

- 1) Выявлена зависимость между коэффициентом невязки и таксономическим положением организмов. У организмов одного рода коэффициенты невязки принимают близкие значения. Корреляция между невязками разных порядков одного отдела выше, чем между невязками разных отделов.
- 2) Значение невязки падает экспоненциально с ростом толщины словаря, для любых организмов.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

1. Усанов Н.Н., Усанов Н.Г. Симметричные структуры в комплементарных цепях геномных ДНК и возможные алгоритмы их организации. Москва : б.н., 2011 г.
2. Капун Е. Д. Разработка метода сравнения нуклеотидных последовательностей путем разбиения на фрагменты. Санкт-Петербург : б.н., 2010 г.
3. Патрушев Л. И., Минкевич И. Г. Проблема размера геномов эукариот. // Успехи биологической химии – 2007 – Том 47, С. 293–370.
4. Koonin E. V., Wolf Y. I. Genomics of bacteria and archaea: the emerging dynamic view of the prokaryotic world (англ.) // Nucleic acids research. — 2008. — Vol. 36, no. 21. — P. 6688-6719.
5. Atsushi Nakabachi, Atsushi Yamashita, Hidehiro Toh, Hajime Ishikawa, Helen E. Dunbar, Nancy A. Moran, Masahira Hattori. The 160-Kilobase Genome of the Bacterial Endosymbiont Carsonella // Science. 2006. V. 314. P. 267.
6. Aleksey Zimin, Kristian A. Stevens, Marc W. Crepeau, Ann Holtz-Morris, Maxim Koriabine, Guillaume Marçais, Daniela Puiu, Michael Roberts, Jill L. Wegrzyn, Pieter J. de Jong, David B. Neale, Steven L. Salzberg, James A. Yorke, Charles H. Langley. Sequencing and Assembly of the 22-Gb Loblolly Pine Genome // GENETICS March 1, 2014 V. 196 no. 3 P. 875-890
7. Бугаенко Н.Н., Горбань А.Н., Садовский М.Г. Определение информационной емкости нуклеотидных последовательностей. Красноярск: б.н., 1997 г.

8. Садовский М.Г. Об информационной емкости символьных последовательностей. // Вычислительные технологии – 2005 - Том 10, № 4 – С. 82- 90.
9. Albrecht-Bühler G. Inversions and inverted transpositions as the basis for an almost universal “format” of genome sequences // Genomics, 2008, vol. 90, pp. 297–305.
10. D. Mitchell, R. Bridge A test of Chargaff's second rule // Biochem. Biophys. Res. Commun., 340 (2006), pp. 90–94
11. Зайцева Н.А. Нарушение второго правила Чаргаффа у митохондриальных геномов и его связь с таксономией носителя. Красноярск, СФУ.
12. Гребнев Я.В., Садовский М.Г. Второе правило чаргаффа и симметрия геномов // Фундаментальные исследования – 2014, № 12 – С. 965-968.